



# An asymptotic preserving scheme for strongly anisotropic elliptic problems

Pierre Degond, Fabrice Deluzet, Claudia Negulescu

## ► To cite this version:

Pierre Degond, Fabrice Deluzet, Claudia Negulescu. An asymptotic preserving scheme for strongly anisotropic elliptic problems. 2009. hal-00371443v2

**HAL Id: hal-00371443**

**<https://hal.science/hal-00371443v2>**

Preprint submitted on 31 Aug 2009

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# An asymptotic preserving scheme for strongly anisotropic elliptic problems

Pierre Degond<sup>†‡</sup>      Fabrice Deluzet<sup>†‡</sup>      Claudia Negulescu<sup>§</sup>

August 31, 2009

## Abstract

In this article we introduce an asymptotic preserving scheme designed to compute the solution of a two dimensional elliptic equation presenting large anisotropies. We focus on an anisotropy aligned with one direction, the dominant part of the elliptic operator being supplemented with Neumann boundary conditions. A new scheme is introduced which allows an accurate resolution of this elliptic equation for an arbitrary anisotropy ratio.

## 1 Introduction

The objective of this paper is to introduce an efficient and accurate numerical scheme to solve a strongly anisotropic elliptic problem of the form

$$\begin{cases} -\nabla \cdot (\mathbb{A} \nabla \phi) = f, & \text{in } \Omega \\ \phi = 0 & \text{on } \partial\Omega_D, \quad \partial_z \phi = 0 & \text{on } \partial\Omega_z, \end{cases} \quad (1)$$

where  $\Omega \subset \mathbb{R}^2$  or  $\Omega \subset \mathbb{R}^3$  is a domain, with boundary  $\partial\Omega = \partial\Omega_D \cup \partial\Omega_z$  and the diffusion matrix  $\mathbb{A}$  is given by

$$\mathbb{A} = \begin{pmatrix} A_\perp & 0 \\ 0 & \frac{1}{\varepsilon} A_z \end{pmatrix}.$$

The terms  $A_\perp$  and  $A_z$  are of the same order of magnitude, whereas the parameter  $0 < \varepsilon < 1$  can be very small, provoking thus the high anisotropy of the problem. In the present paper the considered anisotropy direction is fixed and is aligned with the  $z$ -axis of a Cartesian coordinate system. The method presented here is extended in some forthcoming works to more general anisotropies [9].

Anisotropic problems are common in mathematical modeling and numerical simulation. Indeed they occur in several fields of applications such as flows in porous media [3, 16], semiconductor modelling [24], quasi-neutral plasma simulations [11], image processing [29, 28], atmospheric or oceanic flows [27], and so on, the list being not exhaustive. More specifically high anisotropy aligned with one direction may occur in shell problems or simulation in stretched media. The initial motivation for the present work is closely related to magnetized plasma simulations such as atmospheric [18, 21] or inertial fusion plasmas [7, 12] or plasma thrusters [1]. In this context, the medium is structured by the magnetic field. Indeed, the

---

<sup>†</sup>Université de Toulouse, UPS, INSA, UT1, UTM, Institut de Mathématiques de Toulouse, F-31062 Toulouse, France

<sup>‡</sup>CNRS, Institut de Mathématiques de Toulouse UMR 5219, F-31062 Toulouse, France

<sup>§</sup>CMI/LATP, Université de Provence, 39 rue Frédéric Joliot-Curie 13453 Marseille cedex 13

motion of charged particles in planes perpendicular to the magnetic field is governed by a fast gyration around the magnetic field lines. This explains the large number of collisions the particles encounter in the perpendicular plane, whereas the dynamic in the parallel direction is rather undisturbed. As a consequence the particle mobilities in the perpendicular and parallel directions differ by many orders of magnitude. In the context of ionospheric plasma modelling [6, 17], the ratio of the aligned and transverse mobilities (denoted in this paper by  $\varepsilon^{-1}$ ) can be as huge as ten to the power ten. The relevant boundary conditions in many fields of application are periodic (for instance in simulations of tokamak plasmas on a torus) or Neumann boundary conditions (see for instance [5] for atmospheric plasmas). The system (1) is thus a good model to elaborate a robust numerical method.

The main difficulties with the resolution of problem (1) are of numerical nature, as solving this singular perturbation problem for small  $0 < \varepsilon \ll 1$  is rather delicate. Indeed, replacing in the anisotropic elliptic equation  $\varepsilon$  by zero, yields an ill-posed problem, which has an infinite number of solutions (namely all functions which are constant in the  $z$ -direction). This feature is translated in the discrete case (after the discretization of the problem) into a linear system which is very ill-conditioned for  $\varepsilon \ll 1$ , due to the different order of magnitudes of the various terms. As a consequence standard numerical methods for the resolution of linear systems lead to important numerical costs and unacceptable numerical errors.

More generally, this numerical difficulty arises when the boundary conditions supplied to the dominant  $O(1/\varepsilon)$  operator lead to an ill-posed problem with a multiplicity of solutions. This is the case for Neumann boundary conditions, but also of periodic boundary conditions. If instead, the boundary conditions are such that the dominant operator gives a well-posed problem with a unique solution, this difficulty vanishes as the leading operator alone will suffice to completely determine the limit solution. In this case, one can resort to standard methods. This is the case of Dirichlet or Robin boundary conditions. In spite of the fact that the problem addressed in the present paper arises only with specific boundary conditions, it has a considerable impact in many physics problem, such as plasmas, geophysical flows, plate and shells, etc. In this paper, we will focus on Neumann boundary conditions because they represent a larger range of physical applications, but we could address periodic boundary conditions in a similar way.

Numerical methods for anisotropic elliptic problems have been extensively investigated in the literature. Depending on the underlying physics, distinct numerical methods are developed. For example domain decomposition (Schur complement) and multigrid techniques, using multiple coarse grid corrections are adapted to anisotropic equations in [14, 22] and [13, 25]. For anisotropy aligned with one (or two directions), point (or plane) smoothers are shown to be very efficient [23]. A problem very similar to (1) is addressed in [15], treated via a parametrisation technique, and seems to give good results for rather large anisotropy ratios. However, these techniques are only developed in the context of an elliptic operator with a dominant part supplemented with Dirichlet boundary conditions.

An alternative approach for dealing with highly anisotropic problems is based on a mathematical reformulation of the continuous problem, in order to obtain a more harmless problem, which can be solved numerically in an uncomplicated manner. In this category can be situated for example asymptotic models, describing for small values of the asymptotic parameter  $\varepsilon$  the evolution of an approximation  $\tilde{\phi}$  of the solution of (1) [5, 20]. However, these asymptotic models are precise only for  $\varepsilon \ll 1$ , and cannot be used on the whole range of values covered by the physical parameter  $\varepsilon$ . Thus model coupling methods have to be employed. In sub-domains where the limit model is no longer valid, the original model has to be used, which means that a model coupling strategy has to be developed. However the coupling strategy requires the existence of an area where both models are valid and still demands an accurate numerical method for the resolution of the original model (i.e. the anisotropic elliptic problem) with large anisotropies. This can be rather undesirable.

In this paper, we present an original numerical algorithm belonging to the second approach. A reformulation of the continuous problem (1) will permit us to solve this problem in an inexpensive way and accurately enough, independently of the parameter  $\varepsilon$ . This scheme is related to the *Asymptotic Preserving* numerical method introduced in [19]. These techniques are designed to provide computations in various regimes without any restriction on the discretization meshes and with the additional property to converge towards the solution of the limit problem when the asymptotic parameter goes to zero. The derivation of such Asymptotic Preserving methods requires first the identification of the limit model. For singular perturbation problems, a reformulation of the problem is required in order to derive a set of equations containing both the initial and the limit model with a continuous transition from one regime to another, according to the values of the parameter  $\varepsilon$ . This reformulated system of equations sets the foundation of the AP-scheme. Other singular perturbations have already been explored in previous studies, for instance quasi-neutral or gyro-fluid limits [10, 12]. These techniques have been first introduced for non-stationary systems of equations, for which the time discretization must be studied with care in order to guarantee the asymptotic preserving property. For the anisotropic elliptic equation investigated in this article, we only need to precise the reformulated system and provide a discretization of this one.

The outline of this paper is the following. Section 2 of this article presents first the initial anisotropic elliptic model. In the remainder of this paper, it will be referred to as the *Singular-Perturbation* model (P-model). The reformulated system (referred to as the *Asymptotic Preserving* formulation or AP-formulation) is then derived. It relates on a decomposition of the solution  $\phi(x, z)$  according to its *mean part*  $\bar{\phi}(x)$  along the  $z$  coordinate and a *fluctuation*  $\phi'(x, z)$  consisting of a correction to the mean part needed to recover the full solution. The mean part  $\bar{\phi}(x)$  is solution of an  $\varepsilon$ -independent elliptic problem, and the fluctuation  $\phi'(x, z) = \phi(x, z) - \bar{\phi}(x)$  is given by a well-posed  $\varepsilon$ -dependent elliptic problem. The advantage is that the  $\varepsilon$ -dependent problem for the fluctuation is well-posed and solvable in an inexpensive way, and this uniformly in  $\varepsilon$ . In the limit  $\varepsilon \rightarrow 0$  the AP-formulation reduces to the so called *Limit* model (*L-model*), whose solution is an acceptable approximation of the P-model solution for  $\varepsilon \ll 1$ . The present derivation is carried out in the framework of an anisotropy aligned along one axis of a Cartesian coordinate system. In the context of magnetized plasma simulations, this initial work is extended in a forthcoming work for the three dimensional case in curvilinear coordinates, designed to fit a more complex magnetic field topology (i.e. anisotropy direction) [6]. The main constraints of this method reside in the construction of the mean part which necessitates the integration of the solution along the anisotropy direction. This operation is easily carried out in the context of coordinates adapted with the anisotropy direction. However, an extension of the techniques presented here is currently developed for non-adapted coordinates [9].

Section 3 is devoted to the numerical implementation of the AP-formulation. Numerical results are then presented for a test case, and the three approaches (AP-formulation, straight discretization and resolution of the P-model and L-model) are compared according to the precision of the approximation for different values of  $\varepsilon$ . In section 4 we shall rigorously analyse the convergence of the AP-scheme. Error estimates will be established which underline the advantages of the AP-scheme as compared to the initial Singular Perturbation model and the Limit model.

Current research directions are concerned with the adaptation of the present technique to the case of arbitrary spatially varying anisotropies, without adaptation of the coordinate system to the direction of the anisotropy. These developments will allow the treatment on nonlinear problems, when the diffusion tensor (and its principal directions) depend on the solution itself. This treatment will involve iterative methods which, at each iterate, will

reduce the problem to the solution of a linear anisotropic diffusion problem.

## 2 The asymptotic preserving formulation

For simplicity we shall consider in this paper the two-dimensional problem, posed on a rectangular domain  $\Omega = \Omega_x \times \Omega_z$ , where  $\Omega_x \subset \mathbb{R}$  and  $\Omega_z \subset \mathbb{R}$  are intervals. The ideas exposed here can be extended without any problems to the more physical three-dimensional domain, with two transverse directions  $(x, y)$  and an anisotropy direction aligned with the  $z$ -direction. In this section we introduce the Singular Perturbation Model, the Limit Model and the Asymptotic Preserving formulation.

### 2.1 The Singular Perturbation Model (P-model)

The main concern of this paper is the numerical resolution of the following anisotropic, elliptic problem, called in the sequel Singular Perturbation Model

$$(P) \quad \begin{cases} -\nabla \cdot (\mathbb{A} \nabla \phi) = f, & \text{in } \Omega, \\ \frac{\partial \phi}{\partial z} = 0 & \text{on } \Omega_x \times \partial\Omega_z, \quad \phi = 0 & \text{on } \partial\Omega_x \times \Omega_z. \end{cases} \quad (2)$$

The anisotropy of the media is modeled via the definition of the diffusion matrix  $\mathbb{A}$

$$\mathbb{A} = \begin{pmatrix} A_\perp & 0 \\ 0 & \frac{1}{\varepsilon} A_z \end{pmatrix}, \quad (3)$$

where  $A_\perp(x, z)$  and  $A_z(x, z)$  are given functions with comparable order of magnitudes. The source term  $f(x, z)$  is given and the parameter  $\varepsilon$  is small compared to both  $A_\perp$  as well as  $A_z$ . The medium becomes more anisotropic as the value of  $\varepsilon$  goes to zero.

### 2.2 The limit regime (L-model)

In this section we establish that in the limit  $\varepsilon \rightarrow 0$  the solution of the perturbation model converges towards  $\bar{\phi}$ , solution of the L-model defined by

$$(L) \quad \begin{cases} -\frac{\partial}{\partial x} \left( \bar{A}_\perp \frac{\partial \bar{\phi}}{\partial x} \right) = \bar{f}(x), & \text{in } \Omega_x, \\ \bar{\phi} = 0 & \text{on } \partial\Omega_x, \end{cases} \quad (4)$$

where overlined quantities designate averages over the  $z$ -coordinate :

$$\bar{f}(x) = \frac{1}{|\Omega_z|} \int_{\Omega_z} f(x, z) dz.$$

First we can rewrite the P-model as

$$(P) \quad \begin{cases} -\frac{\partial}{\partial x} \left( A_\perp \frac{\partial \phi}{\partial x} \right) - \frac{1}{\varepsilon} \frac{\partial}{\partial z} \left( A_z \frac{\partial \phi}{\partial z} \right) = f, & \text{in } \Omega, \\ \frac{\partial \phi}{\partial z} = 0 & \text{on } \Omega_x \times \partial\Omega_z, \quad \phi = 0 & \text{on } \partial\Omega_x \times \Omega_z, \end{cases} \quad (5)$$

and integrating along the  $z$ -coordinate gives

$$\frac{\partial}{\partial x} \left( \overline{A_\perp \frac{\partial \phi}{\partial x}} \right) = \bar{f}(x). \quad (6)$$

This equation holds for any  $\varepsilon > 0$ . Now, letting formally  $\varepsilon$  tend to zero in (5) yields the reduced model (R-model)

$$(R) \quad \begin{cases} -\frac{\partial}{\partial z} \left( A_z \frac{\partial \phi}{\partial z} \right) = 0, & \text{in } \Omega, \\ \frac{\partial \phi}{\partial z} = 0 & \text{on } \Omega_x \times \partial\Omega_z, \quad \phi = 0 & \text{on } \partial\Omega_x \times \Omega_z. \end{cases} \quad (7)$$

The functions verifying this ill-posed R-model are constant along the  $z$ -coordinate. Thus including this asymptotic limit property into equation (6) gives rise to the L-model (4), verified by the solution of the Singular Perturbation model in the limit  $\varepsilon \rightarrow 0$ .

**Remark 2.1** *The L-model is the singular limit of the original P-model (2). It provides an accurate approximation of the P-solution only for small values of  $\varepsilon$ . The P-model is valid for all  $0 < \varepsilon < 1$ , but numerically impracticable for  $\varepsilon \ll 1$ . Indeed working with a finite precision, the asymptotic model degenerates into the R-model defined by (7) as  $\varepsilon$  vanishes. This R-model is ill-posed since it exhibits an infinite amount of solutions  $\phi = \tilde{\phi}(x)$ , depending only on the variable  $x$ . This implies that the discretization matrix derived from the P-model is very ill-conditioned for small  $0 < \varepsilon \ll 1$ . This point is addressed by the numerical experiments of section 3.2. Consequently, in a domain where  $\varepsilon$  varies significantly, a model coupling method has to be developed in order to exploit the validity of each model, the P- and L-model. This can be rather undesirable. In the next section we shall present an alternative approach, which is based on a reformulation of the Singular-Perturbation model providing a means of computing an accurate numerical approximation of the solution for all values  $0 < \varepsilon < 1$ .*

**Remark 2.2** *The asymptotics is totally different in the case of Dirichlet boundary conditions. In this case, the R-model is well posed, with a unique solution, and there is no difficulty anymore. Any standard numerical solution of the P-model will converge to that of the R-model. In other words, with Dirichlet boundary conditions, the perturbation becomes regular and the limit solution is fully determined by the formal limit system. The situation and the difficulty addressed in the present paper require that the R-model be ill-posed. This is the case with Neumann boundary conditions (which is the framework chosen here) but also with periodic boundary conditions, or any other boundary condition which would result in an ill-posed R-model.*

### 2.3 The Asymptotic Preserving reformulation (AP-formulation)

In order to circumvent the just described numerical difficulties in handling the Singular Perturbation model, we introduce a reformulation, which permits a transition from the initial P-model to its singular limit (L-model), as  $\varepsilon \rightarrow 0$ .

For this, we shall decompose each quantity  $f(x, z)$  into its mean value  $\bar{f}(x)$  along the  $z$  coordinate and a fluctuation part  $f'(x, z)$ . For simplicity reasons let in the following  $\Omega_x := (0, L_x)$  and  $\Omega_z := (0, L_z)$ . Then

$$f(x, z) = \bar{f}(x) + f'(x, z), \quad (8)$$

with

$$\bar{f}(x) := \frac{1}{L_z} \int_0^{L_z} f(x, z) dz, \quad f'(x, z) := f(x, z) - \bar{f}(x). \quad (9)$$

Note that we have the following properties

$$\bar{f}' = 0, \quad \overline{(\partial f / \partial x)} = \partial \bar{f} / \partial x, \quad \overline{f'g} = \bar{f}'\bar{g} + \overline{f'g'}, \quad (10)$$

$$\partial f / \partial z = \partial f' / \partial z, \quad (\partial f / \partial x)' = \partial f' / \partial x, \quad (fg)' = f'g' - \overline{f'g'} + \bar{f}g' + f'\bar{g}. \quad (11)$$

Taking now the mean of the elliptic equation (5) along the  $z$ -coordinate, we get thanks to (10) and (11), an equation for the evolution of the mean part  $\bar{\phi}(x)$

$$(AP1) \quad \begin{cases} -\frac{\partial}{\partial x} \left( \bar{A}_\perp \frac{\partial \bar{\phi}}{\partial x} \right) = \bar{f} + \frac{\partial}{\partial x} \left( A'_\perp \frac{\partial \bar{\phi}}{\partial x} \right), & \text{in } \Omega_x, \\ \bar{\phi} = 0 & \text{on } \partial\Omega_x. \end{cases} \quad (12)$$

Subtracting from (5) this mean equation (12), gives rise to the evolution equation for the fluctuation part  $\phi'(x, z)$

$$(AP2) \quad \begin{cases} -\frac{\partial}{\partial z} \left( A_z \frac{\partial \phi'}{\partial z} \right) - \varepsilon \frac{\partial}{\partial x} \left( A_\perp \frac{\partial \phi'}{\partial x} \right) + \varepsilon \frac{\partial}{\partial x} \left( A'_\perp \frac{\partial \phi'}{\partial x} \right) = \\ \varepsilon f' + \varepsilon \frac{\partial}{\partial x} \left( A'_\perp \frac{\partial \bar{\phi}}{\partial x} \right), & \text{in } \Omega, \\ \frac{\partial \phi'}{\partial z} = 0 & \text{on } \Omega_x \times \partial\Omega_z, \quad \phi' = 0 & \text{on } \partial\Omega_x \times \Omega_z, \\ \bar{\phi}' = 0, & \text{in } \Omega_x. \end{cases} \quad (13)$$

Thus we have replaced the resolution of the initial Singular Perturbation model (5) by the resolution of the system (12)-(13), which will be done iteratively. Starting from a guess function  $\phi'$ , equation (12) gives the mean value  $\bar{\phi}(x)$ , which inserted in (13) shall give the fluctuation part  $\phi'(x, z)$  and so on.

The constraint  $\bar{\phi}' = 0$  in (13) (which is automatic for  $\varepsilon > 0$ , as explained in Remark 2.3) has the essential consequence that the conditioning of the discretized system becomes  $\varepsilon$ -independent, because the problem (13) reduces in the limit  $\varepsilon \rightarrow 0$  to the system

$$\begin{cases} -\frac{\partial}{\partial z} \left( A_z \frac{\partial \phi'}{\partial z} \right) = 0, & \text{in } \Omega, \\ \frac{\partial \phi'}{\partial z} = 0 & \text{on } \Omega_x \times \partial\Omega_z, \quad \phi' = 0 & \text{on } \partial\Omega_x \times \Omega_z, \\ \bar{\phi}' = 0 & \text{in } \Omega_x, \end{cases} \quad (14)$$

which is uniquely solvable, with the solution  $\phi' \equiv 0$ . Inserting this solution in (12), we conclude that the solution of the AP formulation converges for  $\varepsilon \rightarrow 0$  towards the mean value part  $\bar{\phi}(x)$ , computed thanks to the Limit model

$$(L) \quad \begin{cases} -\frac{\partial}{\partial x} \left( \bar{A}_\perp \frac{\partial \bar{\phi}}{\partial x} \right) = \bar{f}(x), & \text{in } \Omega_x, \\ \bar{\phi} = 0 & \text{on } \partial\Omega_x. \end{cases} \quad (15)$$

The AP reformulation (12)-(13) is equivalent to the Singular Perturbation problem (5) and is therefore valid for all  $0 < \varepsilon < 1$ . This new formulation guarantees that, working with a finite precision arithmetic, the computed solution converges in the limit  $\varepsilon \rightarrow 0$  towards the solution of the limit model (4). This is a huge difference with the original Singular Perturbation model which degenerates into an ill-posed problem. Thus, by using the AP-formulation, we expect the computation of the numerical solution to be accurate, uniformly in  $\varepsilon$ .

For the detailed mathematical proofs, we refer to the next section.

**Remark 2.3** The condition  $\overline{\phi'} = 0$  in (13) holds automatically for  $\varepsilon > 0$ , since the right-hand side has zero average along the  $z$ -coordinate. Indeed, let  $\psi$  be the solution of

$$\begin{cases} -\frac{\partial}{\partial z} \left( A_z \frac{\partial \psi}{\partial z} \right) - \varepsilon \frac{\partial}{\partial x} \left( A_\perp \frac{\partial \psi}{\partial x} \right) + \varepsilon \frac{\partial}{\partial x} \left( \overline{A'_\perp} \frac{\partial \psi}{\partial x} \right) = \varepsilon g', & \text{in } \Omega, \\ \frac{\partial \psi}{\partial z} = 0 & \text{on } \Omega_x \times \partial\Omega_z, \quad \psi = 0 & \text{on } \partial\Omega_x \times \Omega_z, \end{cases} \quad (16)$$

with  $\overline{g'} = 0$ . Taking the average along  $z$ , we get

$$\begin{cases} -\frac{\partial}{\partial x} \left( \bar{A}_\perp \frac{\partial \bar{\psi}}{\partial x} \right) = 0, & \text{in } \Omega_x, \\ \bar{\psi} = 0 & \text{on } \partial\Omega_x, \end{cases}$$

and thus  $\bar{\psi} \equiv 0$ , which is nothing but the constraint added in (13).

The computations of the fluctuating part  $\phi'$  via the equation (13) requires the discretization of an integro-differential operator. This means that the discretization matrix will contain dense blocks. However, using (12) the system (AP2) can be rewritten as

$$(AP2') \quad \begin{cases} -\frac{\partial}{\partial z} \left( A_z \frac{\partial \phi'}{\partial z} \right) - \varepsilon \frac{\partial}{\partial x} \left( A_\perp \frac{\partial \phi'}{\partial x} \right) = \\ \hspace{15em} \varepsilon f + \varepsilon \frac{\partial}{\partial x} \left( A_\perp \frac{\partial \bar{\phi}}{\partial x} \right), & \text{in } \Omega, \\ \frac{\partial \phi'}{\partial z} = 0 & \text{on } \Omega_x \times \partial\Omega_z, \quad \phi' = 0 & \text{on } \partial\Omega_x \times \Omega_z, \\ \bar{\phi'} = 0, & \text{in } \Omega_x. \end{cases} \quad (17)$$

In this expression the right-hand side has no longer zero mean value along the  $z$ -coordinate, but the integro-differential operator has disappeared. The associated discretization matrix is thus sparser than that obtained from the system (12). Systems (12)-(13) and (12)-(17) are equivalent.

## 2.4 Mathematical study of the AP-formulation

We establish in this section the mathematical framework of the AP-formulation (12)-(13) and study its mathematical properties. Let us thus introduce the two Hilbert-spaces

$$\mathcal{V} := \{ \psi(\cdot, \cdot) \in H^1(\Omega) / \psi = 0 \text{ on } \partial\Omega_x \times \Omega_z \}, \quad \mathcal{W} := \{ \psi(\cdot) \in H^1(\Omega_x) / \psi = 0 \text{ on } \partial\Omega_x \},$$

with the corresponding scalar-products

$$(\phi, \psi)_\mathcal{V} := \varepsilon (\partial_x \phi, \partial_x \psi)_{L^2} + (\partial_z \phi, \partial_z \psi)_{L^2}, \quad (\phi, \psi)_\mathcal{W} := (\partial_x \phi, \partial_x \psi)_{L^2}, \quad (18)$$

and the induced norms  $\| \cdot \|_\mathcal{V}$ , respectively  $\| \cdot \|_\mathcal{W}$ . For simplicity reasons, we denote in the sequel the  $L^2$  scalar-product simply by the bracket  $(\cdot, \cdot)$ . Defining the following bilinear



forms

$$\begin{aligned}
a_0(\phi', \psi') &:= \int_0^{L_z} \int_0^{L_x} A_z(x, z) \frac{\partial \phi'}{\partial z}(x, z) \frac{\partial \psi'}{\partial z}(x, z) dx dz, \\
a_1(\phi', \psi') &:= \int_0^{L_z} \int_0^{L_x} A_\perp(x, z) \frac{\partial \phi'}{\partial x}(x, z) \frac{\partial \psi'}{\partial x}(x, z) dx dz, \\
a_2(\bar{\phi}, \bar{\psi}) &:= \int_0^{L_x} \bar{A}_\perp(x) \frac{\partial \bar{\phi}}{\partial x}(x) \frac{\partial \bar{\psi}}{\partial x}(x) dx, \\
c(\phi', \bar{\psi}) &:= \int_0^{L_z} \int_0^{L_x} A'_\perp(x, z) \frac{\partial \phi'}{\partial x}(x, z) \frac{\partial \bar{\psi}}{\partial x}(x) dx dz, \\
d(\phi', \psi') &:= \frac{1}{L_z} \int_0^{L_z} \int_0^{L_x} \int_0^{L_z} A'_\perp(x, z) \frac{\partial \phi'}{\partial x}(x, z) \frac{\partial \psi'}{\partial x}(x, \zeta) dz d\zeta dx, \\
b(\bar{P}, \psi') &:= \int_0^{L_x} \bar{P}(x) \int_0^{L_z} \psi'(x, z) dz dx, \\
a(\phi', \psi') &:= a_0(\phi', \psi') + \varepsilon a_1(\phi', \psi') - \varepsilon d(\phi', \psi'),
\end{aligned} \tag{19}$$

permits to rewrite the AP system (12)-(13) under the weak form

$$(AP) \quad \begin{cases} a_2(\bar{\phi}, \bar{\psi}) = (\bar{f}, \bar{\psi}) - \frac{1}{L_z} c(\phi', \bar{\psi}), & \forall \bar{\psi} \in \mathcal{W}, \\ a(\phi', \psi') + b(\bar{P}, \psi') = \varepsilon(f', \psi') - \varepsilon c(\psi', \bar{\phi}), & \forall \psi' \in \mathcal{V}, \\ b(\bar{Q}, \phi') = 0, & \forall \bar{Q} \in \mathcal{W}, \end{cases} \tag{20}$$

where  $\phi'(x, z) \in \mathcal{V}$ ,  $\bar{\phi}(x) \in \mathcal{W}$  as well as  $\bar{P}(x) \in \mathcal{W}$  are the unknowns and  $\psi' \in \mathcal{V}$ ,  $\bar{\psi} \in \mathcal{W}$  and  $\bar{Q} \in \mathcal{W}$  the test functions. It can be observed that the constraint  $\bar{\phi}' = 0$  was introduced via the Lagrange multiplier  $\bar{P}$ . We will see in the next theorem that the weak formulation (20) is equivalent for  $\varepsilon > 0$  to the system

$$\begin{cases} a_2(\bar{\phi}, \bar{\psi}) = (\bar{f}, \bar{\psi}) - \frac{1}{L_z} c(\phi', \bar{\psi}), & \forall \bar{\psi} \in \mathcal{W}, \\ a(\phi', \psi') = \varepsilon(f', \psi') - \varepsilon c(\psi', \bar{\phi}), & \forall \psi' \in \mathcal{V}, \end{cases} \tag{21}$$

$$\tag{22}$$

where the explicit constraint  $\bar{\phi}' = 0$  does not appear. Let us assume in the sequel

**Hypothesis A** *Let the diffusion functions  $A_\perp \in L^\infty(\Omega)$  and  $A_z \in L^\infty(\Omega)$  satisfy*

$$0 < c_\perp \leq A_\perp(x, z) \leq M_\perp, \quad 0 < c_z \leq A_z(x, z) \leq M_z, \quad f.a.a. (x, z) \in \Omega,$$

*with some positive constants  $c_\perp, c_z, M_\perp, M_z$ . Let moreover  $f \in L^2(\Omega)$ .*

The next theorem analyzes the well-posedness of the AP-formulation.

**Theorem 2.4** *For every  $\varepsilon > 0$  the problem (21)-(22) admits under Hypothesis A a unique solution  $(\phi'_\varepsilon, \bar{\phi}_\varepsilon) \in \mathcal{V} \times \mathcal{W}$ , where  $\phi_\varepsilon := \phi'_\varepsilon + \bar{\phi}_\varepsilon$  is the unique solution of the Singular Perturbation model (5). The function  $\phi'_\varepsilon$  has zero mean value along the  $z$ -coordinate, i.e.  $\bar{\phi}'_\varepsilon = 0$  for every  $\varepsilon > 0$ .*

*Consequently,  $(\phi'_\varepsilon, \bar{\phi}_\varepsilon) \in \mathcal{V} \times \mathcal{W}$  is the unique solution of (21)-(22) if and only if  $(\phi'_\varepsilon, \bar{\phi}_\varepsilon, \bar{P}_\varepsilon) \in \mathcal{V} \times \mathcal{W} \times \mathcal{W}$  is a solution of the AP-formulation (20). In this last case, we have  $\bar{P}_\varepsilon = 0$ .*

*Finally, these solutions satisfy the bounds*

$$\|\phi_\varepsilon\|_{H^1(\Omega)} \leq C\|f\|_{L^2(\Omega)}, \quad \|\phi'_\varepsilon\|_{H^1(\Omega)} \leq C\|f\|_{L^2(\Omega)}, \quad \|\bar{\phi}_\varepsilon\|_{H^1(\Omega_x)} \leq C\|f\|_{L^2(\Omega)},$$

with an  $\varepsilon$ -independent constant  $C > 0$ . In the limit  $\varepsilon \rightarrow 0$  there exist some functions  $(\phi'_0, \bar{\phi}_0) \in \mathcal{V} \times \mathcal{W}$ , such that we have the following weak convergences in  $H^1$

$$\phi'_\varepsilon \rightharpoonup_{\varepsilon \rightarrow 0} \phi'_0 \quad \text{in } H^1(\Omega), \quad \bar{\phi}_\varepsilon \rightharpoonup_{\varepsilon \rightarrow 0} \bar{\phi}_0 \quad \text{in } H^1(\Omega_x),$$

and the strong  $L^2$  convergences

$$\phi'_\varepsilon \rightarrow_{\varepsilon \rightarrow 0} \phi'_0 \quad \text{in } L^2(\Omega), \quad \partial_z \phi'_\varepsilon \rightarrow_{\varepsilon \rightarrow 0} \partial_z \phi'_0 \quad \text{in } L^2(\Omega), \quad \bar{\phi}_\varepsilon \rightarrow_{\varepsilon \rightarrow 0} \bar{\phi}_0 \quad \text{in } L^2(\Omega_x),$$

where  $\phi'_0 \equiv 0$  and  $\bar{\phi}_0$  is the unique solution of the Limit model (4).

**Proof:** The Singular Perturbation model (5) and the Limit model (4) are standard elliptic problems and posses under Hypothesis A (and for every  $\varepsilon > 0$ ) unique solutions  $\phi_\varepsilon \in \mathcal{V}$ , respectively  $\bar{\phi} \in \mathcal{W}$ . It is then a simple consequence of the decomposition (9), that the problem (21)-(22) admits a unique solution  $(\phi'_\varepsilon, \bar{\phi}_\varepsilon) \in \mathcal{V} \times \mathcal{W}$ , where  $\bar{\phi}_\varepsilon(x) := \frac{1}{L_z} \int_0^{L_z} \phi_\varepsilon(x, z) dz$  is the mean and  $\phi'_\varepsilon := \phi_\varepsilon - \bar{\phi}_\varepsilon$  the fluctuation part. Thus we have also  $\bar{\phi}'_\varepsilon = 0$ . This property can also be understood from the fact that the right-hand side of (13), denoted in the sequel by  $g$

$$g(x, z) := f'(x, z) + \frac{\partial}{\partial x} \left( A'_\perp(x, z) \frac{\partial \bar{\phi}}{\partial x}(x) \right),$$

has zero mean value along the  $z$ -coordinate. Indeed, taking in (22) test functions  $\psi'(x) \in \mathcal{V}$  depending only on  $x$ , yields immediately that  $\bar{\phi}'_\varepsilon = 0$  for all  $\varepsilon > 0$ .

Standard stability results for elliptic problems yield now the  $\varepsilon$ -independent estimate for the solution of the Singular Perturbation model (5)

$$\|\phi_\varepsilon\|_{H^1(\Omega)}^2 \leq \|\partial_x \phi_\varepsilon\|_{L^2(\Omega)}^2 + \frac{1}{\varepsilon} \|\partial_z \phi_\varepsilon\|_{L^2(\Omega)}^2 \leq C \|f\|_{L^2(\Omega)}^2,$$

implying that  $\|\bar{\phi}_\varepsilon\|_{H^1(\Omega_x)}^2 \leq C \|f\|_{L^2(\Omega)}^2$  and  $\|\phi'_\varepsilon\|_{H^1(\Omega)}^2 \leq C \|f\|_{L^2(\Omega)}^2$ , with a constant  $C > 0$  independent of  $\varepsilon > 0$ . Thus there exist some functions  $(\phi'_0, \bar{\phi}_0) \in \mathcal{V} \times \mathcal{W}$ , such that, up to a subsequence  $\phi'_\varepsilon \rightharpoonup_{\varepsilon \rightarrow 0} \phi'_0$  in  $H^1(\Omega)$  and  $\bar{\phi}_\varepsilon \rightharpoonup_{\varepsilon \rightarrow 0} \bar{\phi}_0$  in  $H^1(\Omega_x)$ . Hence we have

$$\int_0^{L_x} \int_0^{L_z} \phi'_\varepsilon(x, z) \psi(x, z) dx dz \rightarrow_{\varepsilon \rightarrow 0} \int_0^{L_x} \int_0^{L_z} \phi'_0(x, z) \psi(x, z) dx dz, \quad \forall \psi \in \mathcal{V}.$$

Taking here  $\psi(x) \in \mathcal{V}$  depending only on the  $x$ -coordinate, we observe that the feature  $\bar{\phi}'_\varepsilon \equiv 0$  yields the crucial property of the limit solution  $\bar{\phi}'_0 \equiv 0$ . Passing now to the limit  $\varepsilon \rightarrow 0$  in (22), we get that  $\phi'_0$  is solution of

$$a_0(\phi'_0, \psi') = 0, \quad \forall \psi' \in \mathcal{V}, \quad \text{with } \bar{\phi}'_0 = 0 \quad \text{in } \Omega_x,$$

which is the weak form of (14) and implies  $\phi'_0 \equiv 0$ . Finally, passing to the limit in (21), yields that  $\bar{\phi}_0$  is the unique solution of the Limit model (4). Because of the uniqueness of the limit  $(\phi'_0, \bar{\phi}_0)$ , we deduce that the whole sequence  $(\phi'_\varepsilon, \bar{\phi}_\varepsilon)$  converges weakly towards this limit. To conclude the first part of the proof, we shall show the strong  $L^2$  convergences. For this, taking in (22)  $\phi'_\varepsilon$  as test function and passing to the limit  $\varepsilon \rightarrow 0$ , yields  $\partial_z \phi'_\varepsilon \rightarrow 0$  in  $L^2(\Omega)$ . As  $\phi'_\varepsilon \in \mathcal{V}$  and  $\bar{\phi}'_\varepsilon = 0$ , the Poincaré inequality

$$\|\phi'_\varepsilon\|_{L^2} \leq C \|\partial_z \phi'_\varepsilon\|_{L^2},$$

is valid and implies that  $\phi'_\varepsilon \rightarrow 0$  in  $L^2(\Omega)$ . The convergence  $\bar{\phi}_\varepsilon \rightarrow \bar{\phi}_0$  in  $L^2(\Omega_x)$  is immediate by compactity. It remains finally to prove the equivalence between (20) and (21)-(22). This is immediate. Indeed, if  $(\phi'_\varepsilon, \bar{\phi}_\varepsilon) \in \mathcal{V} \times \mathcal{W}$  is solution of (21)-(22), then  $(\phi'_\varepsilon, \bar{\phi}_\varepsilon, 0)$  is solution

of (20). And if  $(\phi'_\varepsilon, \bar{\phi}_\varepsilon, \bar{P}_\varepsilon) \in \mathcal{V} \times \mathcal{W} \times \mathcal{W}$  satisfies (20), then  $\bar{P}_\varepsilon \equiv 0$  (obvious by taking as test function in (20)  $\psi'(x) \in \mathcal{V}$  depending only on  $x$ ) and  $(\phi'_\varepsilon, \bar{\phi}_\varepsilon)$  solves hence (21)-(22). ■

The subject of the next section will be the numerical resolution of the AP-formulation (12)-(13) (or (20)) and this shall be done iteratively via a fixed-point application. Let us thus introduce here the fixed-point map, construct an iterative sequence and analyze its convergence. In the rest of this section, the parameter  $\varepsilon > 0$  shall be considered as fixed. Due to the fact that the two systems (20) and (21)-(22) are equivalent, we shall concentrate on the simpler one, i.e. (21)-(22). Let us define the Hilbert space

$$\mathcal{U} := \{\psi(\cdot, \cdot) \in \mathcal{V} / \bar{\psi} = 0\},$$

associated with the scalar product

$$(\phi, \psi)_* := \int_0^{L_x} \int_0^{L_z} A_z \partial_z \phi \partial_z \psi dz dx + \varepsilon \int_0^{L_x} \int_0^{L_z} A_\perp \partial_x \phi \partial_x \psi dz dx,$$

which is equivalent to the scalar product  $(\cdot, \cdot)_\mathcal{V}$  on  $\mathcal{V}$ , defined by (18).

The fixed-point map  $T : \mathcal{U} \rightarrow \mathcal{U}$  is defined as follows: With  $\phi' \in \mathcal{U}$  we associate  $\bar{\phi} \in \mathcal{W}$ , solution of (21). Then constructing the right-hand side of (22) via this  $\bar{\phi} \in \mathcal{W}$ , we define  $T(\phi')$  as the corresponding solution of (22). Denoting by  $(\phi'_*, \bar{\phi}_*) \in \mathcal{V} \times \mathcal{W}$  the unique solution of (21)-(22), we remark by Theorem 2.4 that  $\phi'_* \in \mathcal{U}$  and that it is the unique fixed-point of the map  $T$ .

**Theorem 2.5** *Let  $\varepsilon > 0$  be fixed and let  $\phi'_* \in \mathcal{U}$  be the unique fixed-point of the application  $T : \mathcal{U} \rightarrow \mathcal{U}$  constructed as follows*

$$\phi' \in \mathcal{U} \xrightarrow{(21)} \bar{\phi} \in \mathcal{W} \xrightarrow{(22)} T(\phi') \in \mathcal{U}.$$

*Then for every starting point  $\phi'_0 \in \mathcal{U}$ , the sequence  $\phi'_k := T(\phi'_{k-1}) = T^k(\phi'_0)$  converges in  $(\mathcal{U}, \|\cdot\|_*)$ , and consequently also in  $(\mathcal{U}, \|\cdot\|_\mathcal{V})$ , towards the fixed-point  $\phi'_* \in \mathcal{U}$  of  $T$ .*

The proof of this theorem is based on the following

**Lemma 2.6** [8] *Let  $(\mathcal{U}, \|\cdot\|_*)$  be a normed space and  $T : \mathcal{U} \rightarrow \mathcal{U}$  a contractive application, i.e.*

$$\|T(\phi) - T(\psi)\|_* < \|\phi - \psi\|_*, \quad \forall \phi, \psi \in \mathcal{U} \quad \text{with} \quad \phi \neq \psi.$$

*Then the set of fixed-points of  $T$ , denoted by  $FP(T)$ , is identical with the set of accumulation points of the sequences  $\{T^k(\phi)\}_{k \in \mathbb{N}}$ , with  $\phi \in \mathcal{U}$ , set which is denoted by  $AP(T)$ . Moreover, these two spaces contain at most one element.*

**Proof of theorem 2.5 :**

The linear application  $T$  is well-defined. The first step  $\phi' \in \mathcal{U} \xrightarrow{(21)} \bar{\phi} \in \mathcal{W}$  is immediate by the Lax-Milgram theorem. For the second step, we remark that for given  $\bar{\phi} \in \mathcal{W}$  the equation

$$a(\theta, \psi') = \varepsilon(f', \psi') - \varepsilon c(\psi', \bar{\phi}), \quad \forall \psi' \in \mathcal{V}, \quad (23)$$

has a unique solution  $\theta \in \mathcal{U}$ . Indeed, we notice first (by taking test functions only depending on the  $x$ -coordinate) that  $\bar{\theta} = 0$ . This enables us to consider instead of (23), the variational formulation

$$m(\theta, \psi') = \varepsilon(f', \psi') - \varepsilon c(\psi', \bar{\phi}), \quad \forall \psi' \in \mathcal{V}, \quad (24)$$

where the bilinear form  $a(\cdot, \cdot)$ , which is not coercive, was replaced by the coercive bilinear form  $m(\cdot, \cdot)$ , given by

$$m(\theta, \psi') := a(\theta, \psi') + \frac{\varepsilon M_\perp}{L_z} \int_0^{L_x} \left[ \int_0^{L_z} \partial_x \theta(x, z) dz \right] \left[ \int_0^{L_z} \partial_x \psi'(x, z) dz \right] dx. \quad (25)$$

Indeed, due to the property  $\bar{\theta} = 0$ , the two equations (23) and (24) are equivalent and this time  $m(\cdot, \cdot)$  is a continuous, coercive bilinear form, as for all  $\psi' \in \mathcal{V}$  we have

$$m(\psi', \psi') \geq \int_0^{L_x} \int_0^{L_z} A_z |\partial_z \psi'|^2 dz dx + \varepsilon \int_0^{L_x} \int_0^{L_z} A_\perp |\partial_x \psi'|^2 dz dx \geq C \|\psi'\|_{\mathcal{V}}^2.$$

Thus the Lax-Milgram theorem implies the existence and uniqueness of a solution  $\theta \in \mathcal{U}$  of the continuous problem (24) and hence also of problem (23). We have shown by this that  $T$  is a well-defined mapping.

Furthermore we know that  $T$  admits, for fixed  $\varepsilon > 0$ , a unique fixed-point, denoted by  $\phi'_* \in \mathcal{U}$ . Let us now suppose that we have shown that  $T$  is contractive. Then lemma 2.6 implies that  $FP(T) = AP(T) = \{\phi'_*\}$ . Thus choosing an arbitrary starting point  $\phi'_0 \in \mathcal{U}$ , and constructing the sequence  $\phi'_k := T(\phi'_{k-1}) = T^k(\phi'_0)$ , we deduce that this sequence has a unique accumulation point  $\phi'_*$  in  $\mathcal{U}$ . This means that the sequence  $\{\phi'_k\}_{k \in \mathbb{N}}$  converges in  $(\mathcal{U}, \|\cdot\|_*)$  towards  $\phi'_*$ . Due to the fact that  $\|\cdot\|_*$  and  $\|\cdot\|_{\mathcal{V}}$  are equivalent norms, we have also the convergence in  $(\mathcal{U}, \|\cdot\|_{\mathcal{V}})$ .

It remains to show that  $T$  is contractive. For this let  $\phi'_1, \phi'_2 \in \mathcal{U}$  be two given, distinct functions. Denoting by  $\phi' := \phi'_1 - \phi'_2$ ,  $\bar{\phi} := \bar{\phi}_1 - \bar{\phi}_2$  (where  $\bar{\phi}_i \in \mathcal{W}$  are the corresponding solutions of (21)) and  $\theta' := T(\phi'_1) - T(\phi'_2)$ , we have to show that  $\|\theta'\|_* < \|\phi'\|_*$ . First we observe that  $\bar{\phi}$  solves

$$a_2(\bar{\phi}, \bar{\psi}) = -\frac{1}{L_z} c(\phi', \bar{\psi}), \quad \forall \bar{\psi} \in \mathcal{W}, \quad (26)$$

and  $\theta'$  is solution of

$$a(\theta', \psi') = -\varepsilon c(\psi', \bar{\phi}), \quad \forall \psi' \in \mathcal{V}. \quad (27)$$

Taking in (26)  $\bar{\phi}$  as test function, gives rise to

$$\begin{aligned} \int_0^{L_x} \bar{A}_\perp |\partial_x \bar{\phi}(x)|^2 dx &= - \int_0^{L_x} \left[ \frac{1}{L_z} \int_0^{L_z} A'_\perp \partial_x \phi'(x, z) dz \right] \partial_x \bar{\phi}(x) dx \\ &= - \int_0^{L_x} \left[ \frac{1}{L_z} \int_0^{L_z} A_\perp \partial_x \phi'(x, z) dz \right] \partial_x \bar{\phi}(x) dx \\ &\leq \frac{1}{\sqrt{L_z}} \left[ \int_0^{L_x} \int_0^{L_z} A_\perp |\partial_x \phi'|^2 dz dx \right]^{1/2} \left[ \int_0^{L_x} \bar{A}_\perp |\partial_x \bar{\phi}|^2 dx \right]^{1/2}. \end{aligned}$$

Thus

$$\left[ \int_0^{L_x} \bar{A}_\perp |\partial_x \bar{\phi}(x)|^2 dx \right]^{1/2} \leq \frac{1}{\sqrt{L_z}} \left[ \int_0^{L_x} \int_0^{L_z} A_\perp |\partial_x \phi'|^2 dz dx \right]^{1/2}.$$

Equally, taking in (27)  $\theta'$  as test function gives rise to

$$\begin{aligned}
& \int_0^{L_x} \int_0^{L_z} A_z |\partial_z \theta'|^2 dz dx + \varepsilon \int_0^{L_x} \int_0^{L_z} A_\perp |\partial_x \theta'|^2 dz \leq -\varepsilon \int_0^{L_x} \int_0^{L_z} A_\perp \partial_x \bar{\phi} \partial_x \theta' dz dx \\
& \leq \varepsilon \left[ \int_0^{L_x} \int_0^{L_z} A_\perp |\partial_x \bar{\phi}|^2 dz dx \right]^{1/2} \left[ \int_0^{L_x} \int_0^{L_z} A_\perp |\partial_x \theta'|^2 dz dx \right]^{1/2} \\
& \leq \varepsilon \sqrt{L_z} \left[ \int_0^{L_x} \bar{A}_\perp |\partial_x \bar{\phi}|^2 dx \right]^{1/2} \left[ \int_0^{L_x} \int_0^{L_z} A_\perp |\partial_x \theta'|^2 dz dx \right]^{1/2}.
\end{aligned} \tag{28}$$

This last inequality yields

$$\begin{aligned}
\int_0^{L_x} \int_0^{L_z} A_z |\partial_z \theta'|^2 dz dx & + \varepsilon \int_0^{L_x} \int_0^{L_z} A_\perp |\partial_x \theta'|^2 dz dx \leq \varepsilon L_z \int_0^{L_x} \bar{A}_\perp |\partial_x \bar{\phi}|^2 dx \\
& \leq \varepsilon \int_0^{L_x} \int_0^{L_z} A_\perp |\partial_x \phi'|^2 dz dx \\
& < \int_0^{L_x} \int_0^{L_z} A_z |\partial_z \phi'|^2 dz dx + \varepsilon \int_0^{L_x} \int_0^{L_z} A_\perp |\partial_x \phi'|^2 dz dx.
\end{aligned}$$

In this last step we would have the “equality” if and only if  $\int_0^{L_x} \int_0^{L_z} A_z |\partial_z \phi'|^2 dz dx = 0$ . This is however only possible for functions depending exclusively on the  $x$ -coordinate,  $\phi'(x)$ , which is in contradiction with the fact that  $\bar{\phi}' = 0$  and  $\phi' \neq 0$ . Thus we have shown that  $\|T(\phi')\|_* < \|\phi'\|_*$  for  $\phi' \neq 0$ ,  $\phi' \in \mathcal{U}$ , which means that  $T$  is a contractive application on  $(\mathcal{U}, \|\cdot\|_*)$ .  $\blacksquare$

### 3 Numerical discretization and simulation results

This part of the paper is concerned with the numerical discretization of the AP-scheme (12)-(13) and the comparison of the simulation results with those obtained via the Singular Perturbation model (5) and the Limit model (4).

#### 3.1 Discretization

The numerical resolution of the Asymptotic Preserving system (12)-(13) is done by means of the standard finite element method.

Let us recall the variational formulation of the AP-formulation

$$\begin{cases} a_2(\bar{\phi}, \bar{\psi}) = (\bar{f}, \bar{\psi}) - \frac{1}{L_z} c(\phi', \bar{\psi}), & \forall \bar{\psi} \in \mathcal{W}, \\ a(\phi', \psi') + b(\bar{P}, \psi') = \varepsilon(f', \psi') - \varepsilon c(\psi', \bar{\phi}), & \forall \psi' \in \mathcal{V}, \\ b(\bar{Q}, \phi') = 0, & \forall \bar{Q} \in \mathcal{W}, \end{cases} \tag{29}$$

with the notation of section 2. Here  $\phi'(x, z) \in \mathcal{V}$ ,  $\bar{\phi}(x) \in \mathcal{W}$  as well as  $\bar{P}(x) \in \mathcal{W}$  are the unknowns and  $\psi' \in \mathcal{V}$ ,  $\bar{\psi} \in \mathcal{W}$  and  $\bar{Q} \in \mathcal{W}$  the test functions.

The introduction of the Lagrange multiplier  $\bar{P}(x)$  was explained in a simplistic manner in the preceding sections and will be analyzed in more details in section 4. Due to the equivalence of (29) and (21)-(22), one can comment that the introduction of  $\bar{P}(x)$  is superfluous,

but this is not the case for the discretized equations. The property  $\overline{\phi}' = 0$  is indeed automatically fulfilled since the right-hand side of equation (22) has a zero mean value along the  $z$ -coordinate. However the discrete implementation of this quantity introduces round-off errors which probably will destroy the zero mean value property and justify the introduction of the Lagrange multiplier.

For simplicity reasons we omitted here the  $\varepsilon$ -index of the solution  $(\phi'_\varepsilon, \overline{\phi}_\varepsilon)$ , the parameter  $\varepsilon > 0$  being considered as fixed.

To discretize now the system (29) we introduce the grid

$$0 = x_0 \leq \dots \leq x_n \leq \dots \leq x_{N_x+1} = L_x, \quad 0 = z_1 \leq \dots \leq z_k \leq \dots \leq z_{N_z} = L_z$$

and denote the cells by  $I_n := [x_n, x_{n+1}]$  and  $J_k := [z_k, z_{k+1}]$ . The finite dimensional spaces  $\mathcal{V}_h \subset \mathcal{V}$  and  $\mathcal{W}_h \subset \mathcal{W}$  are constructed as usual, by means of the hat functions ( $Q_1$  finite elements)

$$\chi_n(x) := \begin{cases} \frac{x - x_{n-1}}{x_n - x_{n-1}}, & x \in I_{n-1}, \\ \frac{x_{n+1} - x}{x_{n+1} - x_n}, & x \in I_n, \\ 0, & \text{else} \end{cases}, \quad \kappa_k(x) := \begin{cases} \frac{z - z_{k-1}}{z_k - z_{k-1}}, & z \in J_{k-1}, \\ \frac{z_{k+1} - z}{z_{k+1} - z_k}, & z \in J_k, \\ 0, & \text{else} \end{cases}.$$

Thus we are searching for approximations  $\phi'_h \in \mathcal{V}_h$ ,  $\bar{\phi}_h \in \mathcal{W}_h$  and  $\bar{P}_h \in \mathcal{W}_h$ , which can be written under the form

$$\phi'_h(x, z) = \sum_{n=1}^{N_x} \sum_{k=1}^{N_z} \alpha_{nk} \chi_n(x) \kappa_k(z), \quad \bar{\phi}_h(x) = \sum_{n=1}^{N_x} \beta_n \chi_n(x), \quad \bar{P}_h(x) = \sum_{n=1}^{N_x} \gamma_n \chi_n(x).$$

Inserting these decompositions in the variational formulation (29) and taking as test functions the hat-functions  $\chi_n$  and  $\kappa_k$  gives rise to the following linear system to be solved in order to get the unknown coefficients  $\alpha_{nk}$ ,  $\beta_n$  and  $\gamma_n$

$$A_2 \beta = \mathbf{w}, \tag{30}$$

$$\begin{pmatrix} A_0 + \varepsilon(A_1 - D) & B \\ B^t & 0 \end{pmatrix} \begin{pmatrix} \alpha \\ \gamma \end{pmatrix} = \varepsilon \begin{pmatrix} \mathbf{v} \\ 0 \end{pmatrix}, \tag{31}$$

where the matrices  $A_2 \in \mathbb{R}^{N_x \times N_x}$ ,  $A_0, A_1, D \in \mathbb{R}^{N_x N_z \times N_x N_z}$  and  $B \in \mathbb{R}^{N_x N_z \times N_x}$  correspond to the bilinear forms (19) and the right-hand sides are defined by

$$\mathbf{w}_n := (\bar{f}, \chi_n) - \frac{1}{L_z} c(\phi'_h, \chi_n), \quad \mathbf{v}_{nk} := (f', \chi_n \kappa_k) - c(\chi_n \kappa_k, \bar{\phi}_h) = (g, \chi_n \kappa_k),$$

for all  $n = 1, \dots, N_x$ ;  $k = 1, \dots, N_z$  and

$$g(x, z) := f'(x, z) + \frac{\partial}{\partial x} \left( A'_\perp(x, z) \frac{\partial \bar{\phi}}{\partial x}(x) \right). \tag{32}$$

Solving iteratively the linear systems (30)-(31) permits finally to get the unknown function  $\phi_h(x, z) = \bar{\phi}_h(x) + \phi'_h(x, z)$ . The convergence of the iterations was proved for the continuous case in theorem 2.5 and can be identically adapted for the discrete case.

### 3.2 Numerical results

In this section we shall compare the numerical results obtained by the discretization of the Singular Perturbation model, the Limit model and the just presented Asymptotic Preserving

reformulation. With this aim, we consider a test case where the exact solution is known. Let thus

$$\phi_e(x, z) := \sin\left(\frac{2\pi}{L_x}x\right) + \varepsilon \cos\left(\frac{2\pi}{L_z}z\right) \sin\left(\frac{2\pi}{L_x}x\right), \quad (33)$$

be the exact solution of problem (5), where we choose  $A_\perp(x, y) = c_1 + xz^2$  and  $A_z(x, z) = c_2 + xz$ , with two constants  $c_1 > 0$ ,  $c_2 > 0$ . The numerical experiments are performed with  $L_x = L_z = 10$  and  $c_1 = c_2 = L_z$ . The exact right-hand side  $f$  is computed by inserting (33) in (5). We denote by  $\phi_P$ ,  $\phi_L$  and  $\phi_A$ , respectively, the numerical solutions of the Singular Perturbation model (5), the Limit model (4) and the Asymptotic Preserving formulation (12)-(13). The comparison will be done in the  $l^2$ -norm, that means

$$\|\phi_e - \phi_{num}\|_2 = \frac{1}{\sqrt{N}} \left( \sum_{i \in \mathcal{G}} |\phi_e(X_i) - \phi_{num,i}|^2 \right)^{1/2}, \quad (34)$$

where  $\phi_{num}$  stands for one of the numerical solutions and  $\phi_e(X_i)$  is the exact solution evaluated in the grid point  $X_i$ . The index  $i$  covers all possible grid indices, reassembled in the set  $\mathcal{G}$ , and  $N$  is the total number of grid points. The linear systems obtained after the discretization of either the P-model, the L-model or the AP-formulation are solved thanks to the same numerical algorithm (MUMPS [2]). The purpose here is not to design a specific preconditioner for the resolution of these linear systems, but to point out the efficiency of the presently introduced AP-method to deal with a large range of anisotropy ratios.

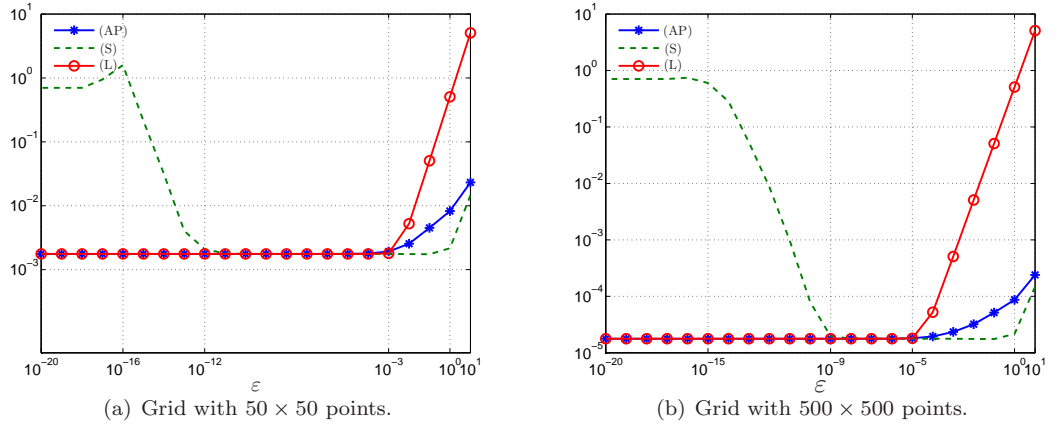


Figure 1: Absolute error in the  $l^2$ -norm between the computed solutions  $\phi_P, \phi_L, \phi_A$  and the exact solution  $\phi_e$ , as a function of  $\varepsilon$  and on different grids. Dashed lines : (S) Standard scheme : discretization of the P-model; Stars : (AP) AP-scheme; Circles : (L) discretization of the L-model.

As can be seen from Table 1 and Figure 1, the finite element resolution of the Singular Perturbation model is precise only for large  $0 < \varepsilon < 1$ , whereas the Limit model is accurate for small  $\varepsilon \ll 1$ . The range of  $\varepsilon$ -values in which both the Singular Perturbation and the Limit models provide an accurate approximation of the solution shrinks as the mesh size is refined. For a coarse grid (with  $50 \times 50$  points see figure 1(a)) this domain ranges from  $10^{-12}$  to  $10^{-3}$  while it is reduced to  $10^{-9} - 10^{-5}$  for the refined  $500 \times 500$  grid (figure 1(b)). This question is determinant for the development of a model coupling strategy. Indeed it requires an intermediate area where both discretized models furnish an accurate approximation and we observe that for refined meshes this area may not exist. This reduction of the validity domain can be explained for both the L-model and P-model but for quite different reasons.

$\varepsilon$	10	1	$10^{-1}$	$10^{-4}$	$10^{-14}$	$10^{-16}$
AP-scheme	$3.4 \cdot 10^{-2}$	$7.8 \cdot 10^{-3}$	$3.8 \cdot 10^{-3}$	$2.7 \cdot 10^{-3}$	$2.7 \cdot 10^{-3}$	$2.7 \cdot 10^{-3}$
S-scheme	$2.8 \cdot 10^{-2}$	$4.5 \cdot 10^{-3}$	$2.8 \cdot 10^{-3}$	$2.7 \cdot 10^{-3}$	$6.6 \cdot 10^{-2}$	1.2
L-model	9.9	$1.0 \cdot 10^1$	$1.0 \cdot 10^{-1}$	$2.8 \cdot 10^{-3}$	$2.7 \cdot 10^{-3}$	$2.7 \cdot 10^{-3}$

Table 1: Absolute error in the  $l^\infty$ -norm for the approximation computed thanks to the AP-scheme, discretized Singular Perturbation and Limit models (S-scheme and L-model) as compared to the exact solution.

The numerical approximation computed via the Limit model is altered by both the discretization error of the numerical scheme and the approximation error introduced by the reduction of the initial Singular Perturbation problem to the Limit problem. For coarse grids, the global error is rapidly dominated by the scheme discretization error, but as the mesh is refined, the approximation error becomes preponderant, as the Limit model is precise only for small  $\varepsilon$ -values. The schemes implemented here are of second order, thus when the mesh size is divided by ten, the discretization error is reduced by one hundred. The global error for the L-model displayed in figure 1(a) does not depend on  $\varepsilon$  as soon as  $\varepsilon < 10^{-3}$ . Below this limit the L-model is able to furnish a better approximation of the solution with vanishing  $\varepsilon$ , however the numerical scheme is not precise enough and consequently the global error does not decrease. For the refined mesh, this discretization error is lowered by two order of magnitudes and the global error is a function of  $\varepsilon$  as long as its value is greater than  $10^{-5}$  (Fig. 1(b)).

The analysis for the Singular Perturbation model is quite complementary. The accuracy of the approximation provided by the P-model is good for large  $\varepsilon$ -values and deteriorates rapidly for small ones. This can be explained by the conditioning of the linear system obtained by the P-model discretization. An estimate of the condition number for the matrix is displayed in figure 2 for two different grid sizes. This conditioning deteriorates with vanishing  $\varepsilon$ -parameter, which is coherent with the fact that, working with a finite-precision arithmetic, the Singular Perturbation model degenerates into an ill-posed problem. This also explains the blow up of the error displayed in figure 1 as soon as the conditioning of the matrix approaches the critical value of the double precision (materialized by the level  $10^{15}$  in Fig. 2). This limit is reached on more refined meshes for larger  $\varepsilon$ -values ( $\varepsilon \approx 10^{-12}$  on a  $50 \times 50$  grid and  $\varepsilon \approx 10^{-10}$  on a  $200 \times 200$  grid). As expected, the P-model, though valid for all  $\varepsilon$ -values, cannot be exploited numerically for small  $\varepsilon$ . The  $\varepsilon$ -region where both the P-model and the L-model are accurate all-together, shrinks dramatically with the size of the mesh, fact which motivates the development of the AP-method.

The condition number estimate of the linear system providing the approximation of the solution for the AP-scheme is also plotted in Figure 2. The conditioning of the system is rather  $\varepsilon$  independent and this is due to the introduction of the Lagrange multiplier, which forces the system in the limit to remain well-posed. The accuracy of the AP-scheme is totally comparable to the P-model for the large values of  $\varepsilon$  and to the L-model for the smallest ones. The AP-formulation is a good tool for computing an approximation for the solution which is accurate uniformly in  $0 < \varepsilon < 1$  and is therefore of great practical interest. Note that this approximation is obtained thanks to an iterative sequence  $\{\phi'_k\}_{k \in \mathbb{N}}$ , constructed with the fixed-point mapping  $T$  defined in theorem 2.5. The convergence of this iterative process is analysed in figure 3 on a  $200 \times 200$  grid for a large value of  $\varepsilon$ . The  $l^2$ -absolute error between the mean respectively the fluctuating parts of the exact solution and the approximation provided by the AP-scheme are plotted as a function of the iteration number. The sequence is initiated with the zero function. With the iterative process, both components converge towards the solution until the precision of the schemes is reached. At this point, after roughly 27 iterations, the approximation can not be improved and a plateau is observed.



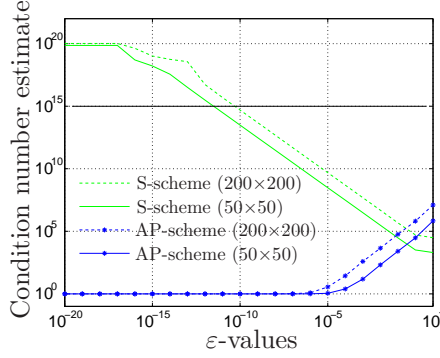


Figure 2: Condition number estimate for the discretization matrices of the Standard (S) and AP schemes (computed by LAPACK [4]) as a function of  $\varepsilon$ . Different grids of  $50 \times 50$  and  $200 \times 200$  points and different  $\varepsilon$ -values are used. Dashed/Plain lines :  $200 \times 200$  /  $50 \times 50$  grid ; Stars : AP-scheme.

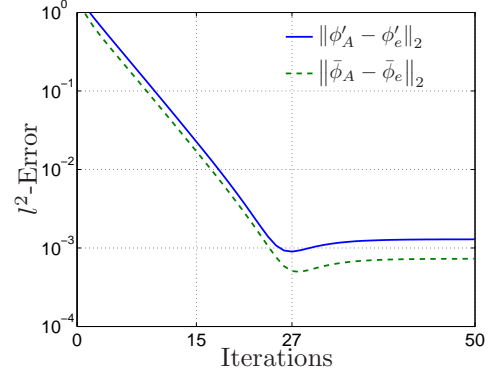


Figure 3: The  $l^2$  absolute error between the exact solution and the numerical approximation computed with the AP-scheme, as a function of the iteration number, with  $\varepsilon = 10$  and a  $200 \times 200$ -mesh. Dashed line : mean part of the solution; Plain line : fluctuating part.

The convergence of this sequence may be improved thanks to classical relaxation techniques. Finally we investigate the positivity of the AP-scheme. With this aim the anisotropic elliptic problem is solved with a positive source term, in this case an approximation of the Dirac  $\delta$ -function. This function denoted  $\delta_a^h$  has a support included in a subset  $([-a, a] \times [-a, a])$ , with  $0 < a < 1$  of the simulation domain  $[-1, 1] \times [-1, 1]$ . Two different parameters  $a$  are chosen,  $a = 10^{-1}$  and  $a = 10^{-2}$ .

The simulation domain is discretized by a  $500 \times 500$  mesh. For the smallest value of  $a$  the support of the function is reduced to 5 cells in each direction. The source term  $\delta_a^h$  is normalized, such that the maximal value of  $\delta_a^h$  grows with vanishing  $a$ -parameter. In table 2 the maxima and minima of the numerical approximations computed by the AP-scheme ( $\phi_A$ ) and the discretized Singular Perturbation model ( $\phi_P$ ) are gathered for the two source functions  $\delta_a^h$ . Only large  $\varepsilon$ -values are considered to verify the positivity of the numerical approximations. Indeed for very small  $\varepsilon$  the solution is reduced to its mean part which is the solution of a classical elliptic problem preserving the maximum principle. This means that the relevant question is related to configurations where the fluctuating part  $\phi'$  has a significant contribution to the elliptic problem solution. In this range of large and intermediate  $\varepsilon$  values, both approximations are comparable. Only slight differences can be observed on the maxima for the smallest  $\varepsilon$ -parameters. The results of table 2 demonstrate the positivity of the approximations computed by either the AP-scheme or the Singular Perturbation model.

	$\varepsilon$	$10^2$	10	1	$10^{-1}$	$10^{-2}$	$10^{-3}$
$a = 10^{-1}$	$\max(\phi_P)$	77.58	3.82	1.63	8.93	7.22	6.93
	$\max(\phi_A)$	77.58	3.82	1.63	8.93	6.89	6.89
	$\min(\phi_P)$	$1.9 \cdot 10^{-7}$	$2.5 \cdot 10^{-7}$	$2.4 \cdot 10^{-2}$	$2.4 \cdot 10^{-2}$	$2.8 \cdot 10^{-2}$	$2.8 \cdot 10^{-2}$
	$\min(\phi_A)$	$1.9 \cdot 10^{-7}$	$2.5 \cdot 10^{-7}$	$2.4 \cdot 10^{-2}$	$2.4 \cdot 10^{-2}$	$2.8 \cdot 10^{-2}$	$2.8 \cdot 10^{-2}$
$a = 10^{-2}$	$\max(\phi_P)$	$1.8 \cdot 10^2$	$7.1 \cdot 10^1$	$2.6 \cdot 10^1$	$1.2 \cdot 10^1$	8.29	7.34
	$\max(\phi_A)$	$1.8 \cdot 10^2$	$7.1 \cdot 10^1$	$2.6 \cdot 10^1$	$1.2 \cdot 10^1$	7.14	7.11
	$\min(\phi_P)$	$1.6 \cdot 10^{-7}$	$2.5 \cdot 10^{-3}$	$2.4 \cdot 10^{-2}$	$2.8 \cdot 10^{-2}$	$2.8 \cdot 10^{-2}$	$2.8 \cdot 10^{-2}$
	$\min(\phi_A)$	$1.6 \cdot 10^{-7}$	$2.5 \cdot 10^{-3}$	$2.4 \cdot 10^{-2}$	$2.8 \cdot 10^{-2}$	$2.8 \cdot 10^{-2}$	$2.8 \cdot 10^{-2}$

Table 2: Maxima and minima of the numerical solutions computed thanks to the AP-scheme ( $\phi_A$ ) and the Singular Perturbation model ( $\phi_P$ ). The elliptic problem is solved with the Dirac  $\delta_a^h$  function as a source term on a  $500 \times 500$  mesh.

## 4 Numerical analysis of the AP-scheme

In this last part of the paper we shall concentrate on the numerical analysis of the  $\mathcal{Q}_1$  finite element scheme introduced in section 3.1 for solving

$$\begin{cases} -\frac{\partial}{\partial z} \left( A_z \frac{\partial \phi}{\partial z} \right) - \varepsilon \frac{\partial}{\partial x} \left( A_\perp \frac{\partial \phi}{\partial x} \right) + \varepsilon \frac{\partial}{\partial x} \left( A'_\perp \frac{\partial \phi}{\partial x} \right) = \varepsilon g, & \text{in } \Omega, \\ \frac{\partial \phi}{\partial z} = 0 & \text{on } \Omega_x \times \partial\Omega_z, \quad \phi = 0 & \text{on } \partial\Omega_x \times \Omega_z, \end{cases} \quad (35)$$

where  $g \in L^2(\Omega)$  is a given function, with mean value along the  $z$ -coordinate equal to zero,  $\bar{g} = 0$ . Moreover we shall explain why we have to introduce the Lagrange multiplier in order to solve numerically this equation. We remark that in contrast to section 3 we omitted for simplicity reasons the primes for  $\phi$ , which indicated the fluctuation functions with zero mean value.

The weak form of (35) is

$$a(\phi, \psi) = \varepsilon(g, \psi), \quad \forall \psi \in \mathcal{V}, \quad (36)$$

or equivalently

$$m(\phi, \psi) = \varepsilon(g, \psi), \quad \forall \psi \in \mathcal{V}, \quad (37)$$

where  $m(\cdot, \cdot)$  is the coercive bilinear form defined in (25). Let us now consider the corresponding discrete problem

$$a(\phi_h, \psi_h) = \varepsilon(g, \psi_h), \quad \forall \psi_h \in \mathcal{V}_h, \quad (38)$$

where the finite dimensional space  $\mathcal{V}_h \subset \mathcal{V}$  was introduced in section 3.1. It can be seen that the property  $\bar{g} = 0$  induces also in the discrete case that  $\overline{\phi_h} = 0$ . Thus, following the same arguments as for the continuous case, we can show that equation (38) is equivalent to

$$m(\phi_h, \psi_h) = \varepsilon(g, \psi_h), \quad \forall \psi_h \in \mathcal{V}_h. \quad (39)$$

The Lax-Milgram theorem implies then the existence and uniqueness of a discrete solution  $\phi_h \in \mathcal{V}_h$ . The next theorem gives an estimate of the discretization error  $\|\phi - \phi_h\|_{\mathcal{V}}$ .

We shall suppose in the sequel, that the diffusion matrices  $A_\perp$ ,  $A_z$  and the function  $f$  are regular enough, to be able to use standard regularity/interpolation results.

**Theorem 4.1** *Let  $\phi \in \mathcal{V}$  be the unique solution of the continuous problem (36) and  $\phi_h \in \mathcal{V}_h$  the unique solution of the discrete problem (38). Both solutions are elements of the normed space  $(\mathcal{U}, \|\cdot\|_{\mathcal{U}})$ , where*

$$\mathcal{U} := \{\psi(\cdot, \cdot) \in \mathcal{V} / \overline{\psi} = 0\} \quad \text{with} \quad \|\psi\|_{\mathcal{U}} := \|\partial_z \psi\|_{L^2(\Omega)}.$$

*Then we have the following discretization error estimate*

$$\|\phi - \phi_h\|_{\mathcal{V}}^2 = \|\partial_z \phi - \partial_z \phi_h\|_{L^2}^2 + \varepsilon \|\partial_x \phi - \partial_x \phi_h\|_{L^2}^2 \leq Ch^2, \quad (40)$$

*with a constant  $C > 0$  independent of  $\varepsilon > 0$ . Moreover, as  $\phi, \phi_h \in \mathcal{U}$ , we have*

$$\|\phi - \phi_h\|_{\mathcal{U}}^2 \leq Ch^2.$$

**Proof:** The fact that both solutions  $\phi$  and  $\phi_h$  belong to the space  $\mathcal{U}$  is an immediate consequence of the fact that the right-hand side of the equation (36) (resp. (38)) satisfies  $\overline{g} = 0$ . The discretization error estimate is rather standard. Denoting by  $\phi_I$  the interpolant of  $\phi$  in the finite dimensional space  $\mathcal{V}_h$ , i.e.

$$\phi_I(x, z) := \sum_{n=1}^{N_x} \sum_{k=1}^{N_z} \phi(x_n, z_k) \chi_n(x) \kappa_k(z),$$

we have due to the coercivity of the bilinear form  $m(\cdot, \cdot)$

$$c \|\phi - \phi_h\|_{\mathcal{V}}^2 \leq m(\phi - \phi_h, \phi - \phi_h) = m(\phi - \phi_h, \phi - \phi_I) \leq c \|\phi - \phi_h\|_{\mathcal{V}} \|\phi - \phi_I\|_{\mathcal{V}}.$$

Thus

$$\|\phi - \phi_h\|_{\mathcal{V}} \leq c \|\phi - \phi_I\|_{\mathcal{V}}.$$

Standard  $\mathcal{Q}_1$  finite element interpolation results [26] yield for the interpolation error

$$\|\partial_x \phi - \partial_x \phi_I\|_{L^2}^2 + \|\partial_z \phi - \partial_z \phi_I\|_{L^2}^2 \leq ch^2 (\|\partial_{xx} \phi\|_{L^2}^2 + \|\partial_{zz} \phi\|_{L^2}^2),$$

and regularity results for the solution  $\phi$  of (36), imply  $\varepsilon^2 \|\partial_{xx} \phi\|_{L^2}^2 + \|\partial_{zz} \phi\|_{L^2}^2 \leq c\varepsilon^2$ . This last estimate can be found by applying standard  $H^2$  regularity results on the solution  $\phi_\varepsilon$  of the initial Singular Perturbation problem (5) (after a change of variable  $\xi := \sqrt{\varepsilon}x$ ) and then exploiting the decomposition  $\phi_\varepsilon = \phi'_\varepsilon + \bar{\phi}_\varepsilon$ . Thus, we have altogether with a constant  $c > 0$  independent of  $\varepsilon > 0$

$$\varepsilon \|\partial_x \phi - \partial_x \phi_h\|_{L^2}^2 + \|\partial_z \phi - \partial_z \phi_h\|_{L^2}^2 \leq ch^2.$$

■

What is important to observe from the error estimate (40) is that for  $\varepsilon \rightarrow 0$  the error  $\|\phi - \phi_h\|_{H^1}$  in the standard  $\varepsilon$ -independent  $H^1$ -norm blows up. This is one argument why the Singular Perturbation model is inaccurate for  $\varepsilon \ll 1$ . However, in the case where  $\phi$  and  $\phi_h$  are elements of the space  $\mathcal{U}$ , we have  $\|\phi - \phi_h\|_{\mathcal{U}} \leq Ch^2$  independently of  $\varepsilon$ , which means that we have convergence of the scheme in  $(\mathcal{U}, \|\cdot\|_{\mathcal{U}})$ , uniformly in  $\varepsilon > 0$ . The Poincaré inequality implies then the uniform convergence in the  $\|\cdot\|_{L^2}$  norm. The AP-scheme is thus equally accurate for every value of  $0 < \varepsilon < 1$ .

The discretization error  $\phi - \phi_h$  is not the only error we are introducing when solving numerically (38) instead of (36). Indeed, (38) is nothing but a linear system

$$M\alpha = v, \quad (41)$$

to be solved to get the unknowns  $\alpha_{nk} := \phi_h(x_n, z_k)$ , where  $v_{nk} := \varepsilon(g, \chi_n \kappa_k)$  and the discrete solution of (38) is then reconstructed as

$$\phi_h(x, z) = \sum_{n=1}^{N_x} \sum_{k=1}^{N_z} \alpha_{nk} \chi_n(x) \kappa_k(z).$$

Unfortunately the implementation of the system (41) introduces round-off as well as approximation errors due for example to the numerical computation of  $a(\chi_n \kappa_k, \chi_r \kappa_p)$ . Thus the numerical resolution of (41) does not yield the exact solution, but an approximation  $(\tilde{\alpha}_{nk})_{nk}$ , solution of the slightly perturbed system

$$M\tilde{\alpha} = \tilde{v}. \quad (42)$$

We are now interested in the error estimate  $\|\phi_h - \tilde{\phi}_h\|_{\mathcal{V}}$ , as a function of the perturbation  $\|v - \tilde{v}\|_2$ , where  $\|\cdot\|_2$  denotes the Euclidean norm in  $\mathbb{R}^{N_x N_z}$ .

**Theorem 4.2** *Let  $\alpha$  be the exact solution of (41) and  $\tilde{\alpha}$  the exact solution of the perturbed system (42). Let  $\phi_h \in \mathcal{V}_h$  and  $\tilde{\phi}_h \in \mathcal{V}_h$  denote the corresponding functions*

$$\phi_h(x, z) = \sum_{n=1}^{N_x} \sum_{k=1}^{N_z} \alpha_{nk} \chi_n(x) \kappa_k(z), \quad \tilde{\phi}_h(x, z) = \sum_{n=1}^{N_x} \sum_{k=1}^{N_z} \tilde{\alpha}_{nk} \chi_n(x) \kappa_k(z).$$

Then we have

$$\varepsilon \|\partial_x \phi_h - \partial_x \tilde{\phi}_h\|_{L^2}^2 + \|\partial_z \phi_h - \partial_z \tilde{\phi}_h\|_{L^2}^2 \leq \frac{c}{\varepsilon} \|v - \tilde{v}\|_2^2, \quad (43)$$

with a constant  $c > 0$  independent of  $\varepsilon > 0$  and  $h > 0$ . However, if both functions  $\phi_h$  and  $\tilde{\phi}_h$  belong to  $\mathcal{U}$ , then we have the  $\varepsilon$ -independent estimate

$$\|\phi_h - \tilde{\phi}_h\|_{\mathcal{U}} \leq c \|v - \tilde{v}\|_2.$$

**Proof:** Let us denote within this proof  $E_{nk} := \alpha_{nk} - \tilde{\alpha}_{nk}$  for  $n = 1, \dots, N_x$ ,  $k = 1, \dots, N_z$  and  $e_h(x, z) := \phi_h(x, z) - \tilde{\phi}_h(x, z)$ , such that

$$e_h(x, z) = \sum_{n=1}^{N_x} \sum_{k=1}^{N_z} E_{nk} \chi_n(x) \kappa_k(z).$$

Moreover let  $N := N_x N_z$  and  $Y \in \mathbb{R}^N$  be an arbitrary vector associated with the function  $y_h(x, z) = \sum_{n=1}^{N_x} \sum_{k=1}^{N_z} Y_{nk} \chi_n(x) \kappa_k(z)$ . Then we have with  $(\cdot, \cdot)_2$  the euclidean scalar product in  $\mathbb{R}^N$  and  $M$  the discretization matrix of (41)

$$\|ME\|_2 = \sup_{Y \in \mathbb{R}^N, Y \neq 0} \frac{(Y, ME)_2}{\|Y\|_2} = \sup_{Y \in \mathbb{R}^N, Y \neq 0} \frac{m(y_h, e_h)}{\|Y\|_2}.$$

Due to the fact that

$$\|Y\|_2 \leq c \|y_h\|_{L^2} \leq \frac{c}{\sqrt{\varepsilon}} \|y_h\|_{\mathcal{V}},$$

we have

$$\|ME\|_2 = \sup_{Y \in \mathbb{R}^N, Y \neq 0} \frac{m(y_h, e_h)}{\|Y\|_2} \geq c\sqrt{\varepsilon} \sup_{y_h \in \mathcal{V}_h, y_h \neq 0} \frac{m(y_h, e_h)}{\|y_h\|_{\mathcal{V}}} \geq c\sqrt{\varepsilon} \|e_h\|_{\mathcal{V}}.$$

Thus we get with a constant  $c > 0$  independent of  $\varepsilon$

$$\|e_h\|_{\mathcal{V}} \leq \frac{c}{\sqrt{\varepsilon}} \|ME\|_2 = \frac{c}{\sqrt{\varepsilon}} \|v - \tilde{v}\|_2.$$

In the case the two functions  $\phi_h$  and  $\tilde{\phi}_h$  belong to  $\mathcal{U}$ , i.e.  $e_h \in \mathcal{U}$ , we can exploit the fact that in  $\mathcal{U}$  the Poincaré inequality gives rise to  $\|Y\|_2 \leq c\|y_h\|_{L^2} \leq c\|y_h\|_{\mathcal{U}}$ . This yields, as  $m(\cdot, \cdot)$  is also coercive on  $\mathcal{U}$ , that

$$\|ME\|_2 = \sup_{Y \in \mathbb{R}^N, Y \neq 0} \frac{m(y_h, e_h)}{\|Y\|_2} \geq c \sup_{y_h \in \mathcal{U}, y_h \neq 0} \frac{m(y_h, e_h)}{\|y_h\|_{\mathcal{U}}} \geq c\|e_h\|_{\mathcal{U}}.$$

and thus the  $\varepsilon$ -independent estimate is proved.  $\blacksquare$

Similarly as for the discretization error, we can deduce from the round-off error estimate (43) that, for  $\varepsilon \rightarrow 0$ , the standard  $H^1$ -norm  $\|\phi_h - \tilde{\phi}_h\|_{H^1}$  explodes. However if we impose that both solutions  $\phi_h$  and  $\tilde{\phi}_h$  are elements of the space  $\mathcal{U}$ , space of functions with mean value along the  $z$ -coordinate equal to zero, then we have the uniform estimate  $\|\phi_h - \tilde{\phi}_h\|_{\mathcal{U}} \leq c\|v - \tilde{v}\|_2$ , and by the Poincaré inequality  $\|\phi_h - \tilde{\phi}_h\|_{L^2} \leq c\|v - \tilde{v}\|_2$ . Unfortunately even if we know that  $\phi_h \in \mathcal{U}$ , this is not necessarily true for  $\tilde{\phi}_h$ , if we discretize (35). But it can be achieved by forcing the numerical solution  $\tilde{\phi}_h$  to satisfy  $\overline{\tilde{\phi}_h} = 0$ . Indeed, this can be done by introducing explicitly in the discrete problem (38) the constraint  $\overline{\tilde{\phi}_h} = 0$ , such that it is much more ingenious to solve instead

$$\begin{cases} a(\phi_h, \psi_h) + b(P_h, \psi_h) = \varepsilon(g, \psi_h), & \forall \psi_h \in \mathcal{V}_h, \\ b(Q_h, \phi_h) = 0, & \forall Q_h \in \mathcal{W}_h, \end{cases} \quad (44)$$

where  $\mathcal{W}_h \subset \mathcal{W}$  was constructed in section 3.1. As mentioned in the continuous case this problem is equivalent for  $\varepsilon > 0$  to the discrete problem (38). If  $\phi_h \in \mathcal{V}_h$  is the unique solution of (38), then  $(\phi_h, 0) \in \mathcal{V}_h \times \mathcal{W}_h$  is a solution of (44). And if  $(\phi_h, P_h) \in \mathcal{V}_h \times \mathcal{W}_h$  solves (44), then  $P_h \equiv 0$  and  $\phi_h \in \mathcal{V}_h$  is the unique solution of (38). This last statement is immediately proved by taking in the variational formulation (44) only  $x$ -dependent test functions  $\psi_h(x) \in \mathcal{V}_h$ . By doing this, we can be sure that the numerical solution  $\tilde{\phi}_h$  of (44) satisfies  $\overline{\tilde{\phi}_h} = 0$ , such that the error  $\|\phi_h - \tilde{\phi}_h\|_{\mathcal{U}}$  is uniformly bounded. This proves that the introduction of the constraint  $\overline{\tilde{\phi}_h} = 0$  in the AP-formulation is crucial and avoids the numerical difficulties associated with the original P-model.

## 5 Conclusion

In this paper we have introduced an Asymptotic Preserving formulation for the resolution of a highly anisotropic elliptic equation. We have shown the advantages of the AP-formulation as compared to the initial Singular Perturbation model and to its limit model, when the asymptotic parameter goes to zero. It came out that the AP-scheme is a powerful tool for the resolution of elliptic problems presenting huge anisotropies along one coordinate, and gives access to the simulation in a very easy and precise manner. The Asymptotic-Preserving method developed here relies on the decomposition of the solution in its mean part along the anisotropy direction, and a fluctuation part. This integration along the anisotropy direction is easily performed in the context of Cartesian coordinate systems with one coordinate aligned with the direction of the anisotropy. In a forthcoming work [9] this procedure is extended to more general anisotropies.

## Acknowledgments

This work has been partially supported by the Marie Curie Actions of the European Commission in the frame of the DEASE project (MEST-CT-2005-021122) and by the CEA-Cesta in the framework of the contracts 'Dynamo-3D' # 4600108543 and 'Magnefig' # 06.31.044.

## References

- [1] J. C. ADAM, J. P. BOEUF, N. DUBUIT, M. DUDECK, L. GARRIGUES, D. GRESILLON, A. HERON, G. HAGELAAR, V. KULAEV, N. LEMOINE, S. MAZOUFFRE, J. PEREZ-LUNA, V. PISAREV, S. TSIKATA, *Physics, simulation, and diagnostics of Hall effect thrusters*, Plasma Phys. Control. Fusion 24, 124041 (2008).
- [2] P. R. AMESTOY, I. S. DUFF, J. KOSTER AND J.-Y. L'EXCELLENT, *A fully asynchronous multifrontal solver using distributed dynamic scheduling*, SIAM Journal of Matrix Analysis and Applications, Vol 23, No 1, pp 15-41 (2001).
- [3] S. F. ASHBY, R. D. FALGOUT, T. W. FOGWELL, A. F. B. TOMPSON, *A numerical solution of groundwater flow and contaminant transport on the CRAY T3D and C90 supercomputers*, Int. J. High Perform. Comp. Appl., Vol. 13 (1999), pp 80-93.
- [4] ANDERSON, E. AND BAI, Z. AND BISCHOF, C. AND BLACKFORD, S. AND DEMMEL, J. AND DONGARRA, J. AND DU CROZ, J. AND GREENBAUM, A. AND HAMMARLING, S. AND MCKENNEY, A. AND SORENSEN, D., *LAPACK Users' Guide, Third Edition*, 1999, pub. Society for Industrial and Applied Mathematics, Philadelphia, PA, ISBN 0-89871-447-8
- [5] C. BESSE, J. CLAUDEL, P. DEGOND, F. DELUZET, G. GALLICE, C. TESSIERAS, *A model Hierarchy for Ionospheric Plasma modeling*, Math. models Methods Appl. Sci, Vol. 14, No. 3 (2004), pp. 393–415.
- [6] C. BESSE, J. CLAUDEL, P. DEGOND, F. DELUZET, G. GALLICE, C. TESSIERAS, *Numerical simulations of the ionospheric striation model in a non-uniform magnetic field*, Comp. Phys. Comm., Vol. 176, No. 2 (2007) pp. 75–90.
- [7] M.A. BEER, S.C. COWLEY, G.W. HAMMETT, *Field-aligned coordinates for nonlinear simulations of tokamak turbulence*, Phys. Plasmas 2 (1995) 2687.
- [8] H. BREZIS, *Points fixes*, Séminaire Choquet, Initiation à l'analyse, Vol. 4 (1964-65), pp. 1–23.
- [9] S. BRULL, P. DEGOND, F. DELUZET, M.-H. VIGNAL, *An asymptotic preserving scheme in the drift limit for the Euler-Lorentz system for a variable magnetic field*, in preparation.
- [10] P. CRISPEL, P. DEGOND AND M.-H VIGNAL, *An asymptotic preserving scheme for the two-fluid Euler-Poisson model in the quasineutral limit*, J. Comput. Phys, 223, (2007) pp. 208–234.
- [11] P. DEGOND, F. DELUZET, L. NAVORET, A.-B. SUN, M.-H. VIGNAL, *Asymptotic-Preserving Particle-In-Cell method for the Vlasov-Poisson system near quasineutrality*, to appear in J. Comput. Phys.
- [12] P. DEGOND, F. DELUZET, A. SANGAM, M.-H. VIGNAL, *An Asymptotic Preserving Scheme for the Euler equations in a strong magnetic field*, J. Comput. Phys, 228, (2009) pp. 3540–3558.
- [13] M.W. GEE, J.J. HU, R.S. TUMINARO, *A new smoothed aggregation multigrid method for anisotropic problems*, Numer. Lin. Alg. with Appl. 16 (2009), pp 19–37.
- [14] L. GIRAUD , AND R.S. TUMINARO, *Schur complement preconditioners for anisotropic problems*, J. Numer. Anal. 19 (1998), pp 1–18.

- [15] P. GUILLAUME, AND V. LATOCHA, *Numerical Convergence of a Parameterisation Method for the Solution of a Highly Anisotropic Two-Dimensional Elliptic Problem* J. Sci. Comput. Vol. 25, 3 (Dec. 2005), pp. 423-444.
- [16] T. Y. HOU, X. H. WU, *A Multiscale Finite Element Method for Elliptic Problems in Composite Materials and Porous Media*, J. Comput Phys. , 134 , 169-189, (1997).
- [17] D. L. HYSSELL, *An overview and synthesis of plasma irregularities in equatorial spread F*, J. Atmos. Solar-Terr. Phys., Vol 62, (2000), pp. 1037–1056.
- [18] M. C. KELLEY, W. E. SWARTZ, J.J. MAKELA, *Mid-Latitude ionospheric fluctuation spectra due to secondary  $E \times B$  instabilities*, J. Atmos. Solar-Terr. Phys., Vol. 66 (2004), pp 1559–1565.
- [19] S. JIN, *Efficient Asymptotic-Preserving (AP) schemes for some multiscale kinetic equations*, SIAM J. Sci. Comp. 21 (1999), pp. 441–454.
- [20] M. J. KESKINEN, *Nonlinear theory of the  $E \times B$  instability with an inhomogeneous electric field*, J. Geophys. Res., Vol. 89, (1984), pp. 3913–3920.
- [21] M. J. KESKINEN, S. L. OSSAKOW, B. G. FEJER, *Three-dimensional nonlinear evolution of equatorial ionospheric spread-F bubbles*, Geophys. Res. Lett., Vol. 30 (2003), pp. 4-1–4-4.
- [22] B.N. KHOROMSKIJ, G. WITTUM, *Robust Schur complement method for strongly anisotropic elliptic equations*, Numer. Linear Algebra Appl., 6, (1999), pp 621–653.
- [23] I.M. LLORENTE AND N.D. MELSON, *Robust multigrid smoothers for three dimensional elliptic equations with strong anisotropies*, ICASE Technical Report: TR-98-37, 1998.
- [24] T. MANKU, A. NATHAN, *Electrical properties of silicon under nonuniform stress*, J. Appl. Phys. **74** (1993), p. 1832.
- [25] Y. NOTAY, *An aggregation-based algebraic multigrid method*, Report GANMN 08-02, Université Libre de Bruxelles, Brussels, Belgium, 2008.
- [26] P.-A. RAVIART, J.-M. THOMAS, *Introduction à l'analyse numérique des équations aux dérivées partielles*, Dunod, Paris, 1998.
- [27] A.-M. TRÉGUIER, *Modélisation numérique pour l'océanographie physique*, Ann. math. Blaise Pascal, tome 9 (2002), no. 2, pp. 345-361.
- [28] W.-W. WANG, X.-C. FENG, *Anisotropic diffusion with nonlinear structure tensor*, Multiscale Model Sim., Vol. 7 (2008), no. 2, pp 963–977.
- [29] J. WEICKERT, *Anisotropic Diffusion in Image Processing* , Teubner, Stuttgart, 1998.